Output files from different CRISPR-Cas detection tools

CASC

• results.txt summary

#seq_id array_id			spacers array_start	array_stop		mean_spacer_len stddev_spacer_len	bonafide	code
1	1	18	2293915 2295042 32.056	0.236	true	7		
1	2	28	2304214 2305952 32.071	0.262	true	7		

Coding scheme (last column):

Code	(Binary)	Cas Protein Hit	Matches Known Repeat	Proper Statistics
0	000	no	no	no
1	001	no	no	Yes
2	010	no	Yes	no
3	011	no	Yes	Yes
4	100	Yes	no	no
5	101	Yes	no	Yes
6	110	Yes	Yes	no
7	111	Yes	Yes	Yes

 bonafide.repeats.fasta bonafide.spacers.fasta non-bonafide.repeats.fasta non-bonafide.spacers.fasta

>1-repeat-1-1
GTGTTCCCCACGGGTGTGGGGATGAACCA
>1-repeat-1-2
GTGTTCCCCACGGGTGTGGGGATGAACCG
>1-repeat-1-3
GTGTTCCCCACGGGTGTGGGGATGAACCG
>1-repeat-1-4
GTGTTCCCCACGGGTGTGGGGATGAACCG
>1-repeat-1-5
GTGTTCCCCACGGGTGTGGGGATGAACCG
>1-repeat-1-6
GTGTTCCCCACGGGTGTGGGGATGAACCG

>1-spacer-1-1
GCAGGATTTGGAGTCGGAGCGTTTACCTACAT
>1-spacer-1-2
GTGAACGATGCGAATCGGGCGGGGGTTTGGTC
>1-spacer-1-3
GCTTGGTAGATGCCTGGAACAAGTCGGTCAGC
>1-spacer-1-4
CGCGCTGTCCTCGGCCTGGACTGGACCATCGA
>1-spacer-1-5
TGGGCGAACGACGGCGCAGCGCA
>1-spacer-1-6
ACAGGGACGTTCTCCGCGAGCTGGCCATCGAC

report.md

```
### Input Summary

- Input File: /mnt/blastdb/emma/crispr_annotation/assembly/unicycler/PaLo1/assembly.fasta
- Number of Sequences = 2
- Number of Bases = 6600643
- Bases per Sequence = 3300322
- Mode = Liberal

### CRISPR Identification
- Putative CRISPR arrays found = 1
- Bona fide CRISPR arrays = 2 (200%)
- Arrays with Cas protein upstream = 1
- Arrays with repeats matching known CRISPR repeats = 1
- Arrays with proper statistics (liberal mode only) = 2
```

casc.log

Commandline: ./casc -i /mnt/blastdb/emma/crispr_annotation/assembly/unicycler/PaLo1/assembly.fasta -o /mnt/

```
Output Directory: /mnt/blastdb/emma/crispr_tools/Try/CASC/PaLo1
```

Mode: Liberal

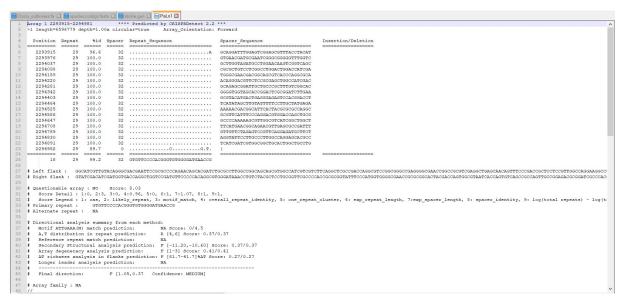
==== CASC Started. Log can be found here: /mnt/blastdb/emma/crispr_tools/Try/CASC/PaLo1/casc.log

==== CASC Finished!

CASC log can be found here: /mnt/blastdb/emma/crispr_tools/Try/CASC/PaLo1/casc.log

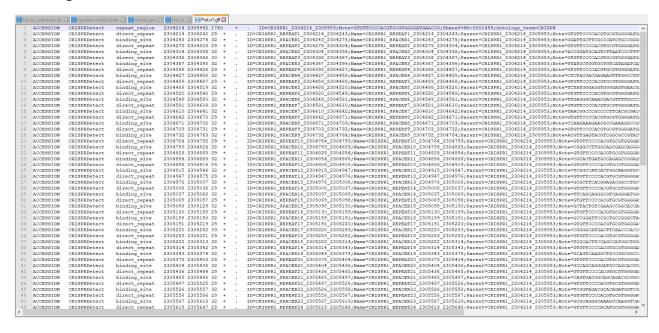
CRISPRDetect

• .txt summary



(but normally CAS data should be there: did not work even though a gbk file was supplied)

- .fp filtered repeats (didn't meet the minimum evidence level): mostly empty
- gff annotation → visualization in CRISPRStudio



CRISPRCasFinder

• .ison -> visualization in CRISPRCasViewer

```
||{
| "Date":"25/10/2019_8:9:22",
"Command": "perl CRISPRCasFinder.pl -levelMin 3 -cas -in /home/Pseudomonas_aeruginosa/PaL
"Sequences":
111
"Id":"1",
"Description": "Unknown",
"AT":33.7792277109783,
"Length": 6596779,
"Summary_CRISPR-Cas":"",
1{
"Name": "1_1",
"Start": 2293915,
"End": 2294980,
"DR Consensus": "GTGTTCCCCACGGGTGTGGGGATGAACCG",
"Repeat_ID": "Unknown",
"DR_Length": 29,
"Spacers": 17,
"Potential Orientation": "+",
"CRISPRDirection": "ND",
"Evidence Level": 4,
"Conservation DRs": 95.7304354599955,
"Conservation_Spacers": 0,
"Regions": [
"Type": "LeftFLANK",
"Start": 2293815,
"End": 2293914,
"Leader": 1,
```

• TSV folder: .tsv+.xls for:

Cas_REPORT

```
1 (Unknown)
### System: CAS-TypeIE (UserReplicon_CAS-TypeIE_1)
#SequenceID Cas-type/subtype Gene status System Type
                                                       Begin End Strand Other information
PaLo1_02144 Cas3_0_I
                      accessory CAS
PaLo1_02145 Cse1_0_IE
                      mandatory
                                  CAS-TypeIE
PaLo1 02146 Cse2 0 IE
                                  CAS-TypeIE
                      mandatory
PaLo1_02147 Cas7_0_IE
                      mandatory
                                 CAS-TypeIE
PaLo1_02148 Cas5_0_IE
                      mandatory
                                  CAS-TypeIE
                                  CAS-TypeIE
PaLo1 02149 Cas6 0 IE
                      mandatory
PaLo1_02150 Cas1_0_IE
                      mandatory
                                  CAS-TypeIE
PaLo1_02151 Cas2_0_IE
                      mandatory
                                  CAS-TypeIE
####Summary system CAS-TypeIE:begin=;end=:{sequenceID=1} : [Cas3_0_I (,,); Cse1_0_IE (,,); Cse2_0 IE
***********************************
### System: CAS (CAS putative)
#SequenceID Cas-type/subtype
                             Gene status System Type
                                                        Begin End Strand Other_information
PaLo1_00479 Cas3_0_I accessory
PaLo1_01544 Cas3_0_I accessory
                                  CAS
                                  CAS
PaLo1_01750 Cas3_0_I
PaLo1_02144 Cas3_0_I
                      accessory
                                  CAS
PaLo1 03983 Cas3 0 I
                                  CAS
                      accessory
PaLo1_04151 Cas3_0_I
                                  CAS
PaLo1_04355 Cas3_1_I
                      accessory
                                  CAS
####$ummary system CAS:begin=;end=:{sequenceID=2} : [Cas3 0 I (,,); Cas3 0 I (,,); Cas3 0 I (,,); Cas3 0 I (,,);
```

Crisprs_REPORT

```
| Sequence | Sequence
```

CRISPR-Cas_summary (CRISPRs and Cas per sequence)

Sequence(s) CRISPR array(s) Nb CRISPRs Evidence-levels Cas cluster(s) Nb Cas Cas Types/Subtypes

1 1_[12293915;2234980] (evidence-level=4), 1_2(2304214;2305952) (evidence-level=4), 1_3(6111088;6111108] (evidence-level=1), 1_4(6516885;6517044) (evidence-level=1), 4 Nb_arrays_evidence-level=20, Nb_arrays_evidence-level=3=0, Nb_arrays_evidence-level=4=0 CAS[;], CAS-TypeIE[;], 2 CAS (n=1), CAS-TypeIE (n=1),

GFF folder: .gff and annotation.gff per sequence in input file

rawCas.fna (not existing if no Cas found)

>1|PaLo1 02144|CAS-TypeIE|Cas3 0 I , GTGTCCGTGGAACTTTGGCAGCAGTGCGTGGATCTTCTCCGCGATGAGCTGCCGTCCCAACAATTCAACACTGGATCCG aatacctcggtcggcttctggaactgctcggtgaacgcggggggtcagttgcccgcgctttccttattaataggcagc AAGCGTAGCCGTACGCCGCGCGCCCATCGTCCCATCGCAGACCCACGTGGCTCCCCGCCTCCGGTTGCTCCGCCCC GGCGCCAGTGCAGCCGGTATCGGCCGCGCCCGTGGTGGTGCCACGTGAAGAGCTGCCGCCAGTGACGACGGCTCCCAGCG GTACGCACCGAGCGCAACGTCCAGGTCGAAGGCGCGCTGAAGCACCAGCTATCTCAACCGTACCTTCACCTTCGAGAA CTTCGTCGAGGGCAAGTCCAACCAGTTGGCCGTGCCGCCTGGCAGGTGGCGGACAACCTCAAGCACGGCTACAACC CGAATTCAAGCGCTTCTACCGCTCGGTGGACGCACTGTTGATCGACGACATCCAGTTCTTCGCCCGTAAGGAGCGCTCCC AGGAGGAGTTCTTCCACACCTTCAACGCCCTTCTCGAAGGCGGCCAGCAGGTGATCCTCACCAGCGACCGCTATCCGAAG GAAATCGAAGGCCTGGAAGAGCGGCTGAAGTCCGCCTTCGGCTGGGCCTGACGGTCGAGCCGTCGAGCCGCAACTGGA AACCCGGGTGGCGATCCTGATGAAGAAGGCCGAGCAGGCGAAGATCGAGCTGCCGCACGACGCGGCCTTCTTCATCGCCC ATCACCATCGAGCTGATTCGCGAGTCGCTGAAGGACTTGTTGGCCCTTCAGGACAAGCTGGTCAGCATCGACAACATCCA GCGCACCGTCGCCGAGTACTACAAGATCAAGATATCCGATCTGTTGTCCAAGCGGCGTTCGCGTTCGGTGGCGCCCGC CACACCACGGTGTTGCACGCCTGTCGTAAGATCGCTCAACTTAGGGAATCCGACGCGGATATCCGCCGAGGACTACAAGAA $\tt CCTGCTGCGTACCCTGACACCTGACGCACGCCCACGAGGCAAGGGACTAGACCATTCACCATTCAACGCGAAGCC$ CTGTTGAAACCGCTGCAACTGGTCGCCGGCGTCGTGGAACGCCGCCAGACATTGCCGGTTCTCCCAACGTCCTGCTGGT

CRISPRs.fna

>1 1 2293915,2294980

>1 2 2304214,2305952

More possible using command-line options (clustering etc.)